

1 UNITED STATES DISTRICT COURT

2 NORTHERN DISTRICT OF CALIFORNIA

3 Before The Honorable Susan van Keulen, Magistrate Judge

4
5 IN RE GOOGLE GENERATIVE)
6 AI COPYRIGHT LITIGATION) No. C 23-03440-EKL
7)

8 San Jose, California
Wednesday, June 18, 2025

9 TRANSCRIPT OF PROCEEDINGS OF THE OFFICIAL ELECTRONIC SOUND
10 RECORDING 10:05 - 11:26/11:43 - 11:49 = 1 HOUR 27 MINUTES

11 APPEARANCES:

12 For Plaintiffs

13 Joseph Saveri Law Firm, LLP
14 601 California Street
Suite 1000
San Francisco, California
94108

15 BY: CHRISTOPHER K.L. YOUNG, ESQ.

16 Bleichmar Fonti & Auld, LLP
17 75 Virginia Road
2nd Floor
White Plains, New York 10603

18 BY: GREGORY MULLENS, ESQ.

19 Blechmar Fonti & Auld, LLP
20 1330 Broadway
Suite 630
Oakland, California 94612

21 BY: ANNE K. DAVIS, ESQ.
22 LESLEY E. WEAVER, ESQ.

23
24 (APPEARANCES CONTINUED ON NEXT PAGE)
25

1 For Defendants:

2 Wilson Sonsini Goodrich &
3 Rosati
4 95 S state Street
Suite 1000
Salt Lake City, Utah 84111
BY: PAUL SAMPSON, ESQ.

5 Wilson Sonsini Goodrich &
6 Rosati
7 650 Page Mill Road
Palo Alto, California 94304
BY: MAURA L. REES, ESQ.
QIFAN HUANG, ESQ.

8 Wilson Sonsini Goodrich &
9 Rosati
10 701 Fifth Avenue
Suite 5100
11 Seattle, Washington 98104
BY: ERIC P. TUTTLE, ESQ.

12 Transcribed by:

Echo Reporting, Inc.
Contracted Court Reporter/
Transcriber
echoreporting@yahoo.com

13
14
15
16
17
18
19
20
21
22
23
24
25

1 Wednesday, June 18, 2025 10:05 a.m.

2 P-R-O-C-E-E-D-I-N-G-S

3 --oOo--

4 THE CLERK: Now calling case number 23CV03440EKL,
5 case In re Google -- I'm sorry, In regarding Google
6 Generative AI Copyright Litigation.

7 If Counsel can state their appearances, starting with
8 the Plaintiffs.

9 MS. WEAVER: Good morning, your Honor. Lesley
10 Weaver, Bleichmar Fonti, for the Plaintiffs.

11 THE COURT: Thank you, Ms. Weaver. Good morning.

12 MS. WEAVER: Good morning.

13 MR. MULLENS: Good morning, your Honor. Greg
14 Mullens, from Bleichmar Fonti, here for the Plaintiffs.

15 THE COURT: Mr. Mullens, hello.

16 MS. DAVIS: Good morning, your Honor. Anne Davis,
17 also from Bleichmar Fonti, for Plaintiffs.

18 THE COURT: Welcome. Ms Davis.

19 MR. YOUNG: Good morning, your Honor. Christopher
20 Young, Joseph Saveri Law Firm, for the Plaintiffs.

21 THE COURT: Thank you, Mr. Young. Welcome.

22 MR. TUTTLE: Good morning, your Honor. Eric
23 Tuttle of Wilson Sonsini, here for the Defendants.

24 THE COURT: Mr. Tuttle, hello.

25 MS. REES: Good morning. Maura Rees from Wilson

1 Sonsini for Defendants.

2 THE COURT: Ms. Rees, good morning.

3 MR. SAMPSON: Good morning, your Honor. Paul
4 Sampson from Wilson Sonsini for the Defendants.

5 THE COURT: Thank you, Mr. Sampson.

6 MR. HUANG: Good morning, your Honor. Qifan Huang
7 from Wilson Sonsini, also for Defendants.

8 THE COURT: Excellent. Mr. Huang, good morning.

9 All right. Thank you all for that. I think I've got
10 everybody checked in. As I mentioned to my CRD, I'm sure
11 you were thrilled to get to use your business cards today.
12 I don't think anybody ever gets to use those anymore. Let
13 me just -- give me just a moment. I just want my law clerks
14 to know I don't have Teams, unfortunately. Whatever is on
15 here is an old version, so I'm just going to open Outlook,
16 if they need me to -- so we'll do it the old fashioned way.
17 If you need to communicate with me to my clerks, you can do
18 so via Outlook. I'm -- probably will recess at some point,
19 so --

20 All right. We have a fair amount of ground to cover,
21 so we're going to start, although Ms. Weaver is anxious to
22 bring me up to date on something.

23 MS. WEAVER: I did just want to -- good morning,
24 your Honor. Lesley Weaver, on behalf of the Plaintiffs. I
25 wanted to let you know that we took your -- the spirit of

1 your order to heart, and we got a counteroffer on the
2 custodian issue from the Defendants yesterday, and we all
3 arrived here early to have further discussions.

4 THE COURT: Fantastic.

5 MS. WEAVER: And Mr. Mullens -- my firm will
6 address that --

7 THE COURT: Okay.

8 MS. WEAVER: -- for our side, and Mr. Young will
9 address source code, depending on what your Honor wants to
10 hear and in what order.

11 THE COURT: Okay.

12 MS. WEAVER: Okay.

13 THE COURT: All right. Well, I don't want to bury
14 the lead, so why don't we go ahead and start with
15 custodians, and let's just see where we are on that. I was
16 going to start with source code, but let's start with
17 custodian.

18 Mr. Mullens?

19 MR. MULLENS: Good morning again, and Greg Mullens
20 for Plaintiffs. Judge, as my colleague just mentioned, this
21 morning, we met and conferred with the Defendants. And let
22 me just get right to it and give you some numbers --

23 THE COURT: Okay.

24 MR. MULLENS: -- so you have them for custodians.
25 Google's, sort of, ask initially was 13 custodians, Judge,

1 and yesterday they countered with also adding seven to that
2 number --

3 THE COURT: Okay.

4 MR. MULLENS: -- so that would be 20.

5 THE COURT: Okay.

6 MR. MULLENS: You know, it took us time to
7 consider that. We had a fair ask, a large ask in terms of
8 our custodians. We certainly submit it was warranted,
9 considering the size and the scope of the case, but we
10 wanted to substantially compromise, as the Court ordered.
11 Upon reflection and consideration and speaking with Defense
12 counsel this morning, we have, from our ask of custodians,
13 cut 20 of our custodians. We now stand at the 13 Google
14 custodians, plus 20 more, which would be 33.

15 And in light of the two declarations that were
16 submitted in support of the source code briefing, Judge, and
17 those are two individuals who are never identified as a
18 Google custodian, clearly very relevant to this case, we're
19 also asking that those two individuals be added. So our
20 number, just to give you a number, is 35.

21 I will say it took a lot for us to reduce our list down
22 by 20 custodians. We're looking to -- or we are looking to
23 stay there as our number. We want to engage with Google if
24 they want to on it. We've done a lot of work in terms of
25 our investigation into these custodians, also using Google's

1 documents which they produced to substantiate our ask. Your
2 Honor accepted and considered for this dispute the
3 organizational charts, which were actually very helpful in
4 identifying unique custodians. And that's all to say that's
5 where we stand at this very moment in terms of the numbers.

6 THE COURT: Okay.

7 MR. MULLENS: And we're looking to certainly
8 answer any questions from the Court about the basis for our
9 custodians or continue speaking with Google if that's what
10 the Court wants.

11 THE COURT: Okay. Mr. Mullens, the seven that --
12 custodians that Google added in their most recent offer,
13 were they seven off of your list?

14 MR. MULLENS: They -- it was unclear, and the
15 answer is no, actually. They said they would offer up
16 seven. They didn't give the names of who. They said they
17 wanted to choose four, and we would choose three. We
18 pressed them on who were at least the four that you were
19 selecting -- from our list. It would be from our list of
20 custodians.

21 THE COURT: I see.

22 MR. MULLENS: But they wouldn't give us the names.
23 And at this juncture, part of our ask for our custodial ask
24 is that we be the ones who pick the custodians and not sort
25 of have a blind decision being made by Google as to who is

1 relevant based upon what we've done here in our compromises.
2 So, right now, it's seven from our list of custodians from
3 Google, but no indication who they are.

4 THE COURT: Okay. And where are the parties --
5 have the parties discussed search terms, or is there a plan
6 once custodians are identified? Remind me where you are in
7 the process.

8 MR. MULLENS: Yeah, search term negotiations are
9 over for both sides.

10 THE COURT: Okay.

11 MR. MULLENS: So we're done with that.

12 THE COURT: Okay. So you have the terms in place.

13 MR. MULLENS: We do, Judge.

14 THE COURT: You're just trying to figure out --

15 MR. MULLENS: Yeah.

16 THE COURT: -- the pool. Well, that's a relief.
17 Okay.

18 And I did look at the chart. Google DeepMind. I'm
19 sure this was explained to me somewhere, but I had a lot of
20 material to review. What does that refer to?

21 And don't worry, Google. I will hear from you too. I
22 know why I'm asking the Plaintiffs.

23 MR. MULLENS: It really is the organization within
24 Google that is responsible for generative AI --

25 THE COURT: Okay.

1 MR. MULLENS: -- development training
2 (indiscernible), but that's --

3 THE COURT: So that's the umbrella, if you will.

4 MR. MULLENS: Yeah, that's the -- sort of the core
5 entity or unit that handles it.

6 THE COURT: Okay. And then the products, some of
7 which are still disputed in the motion to dismiss, those are
8 all under or in the DeepMind --

9 MR. MULLENS: We would say they are, yes, Judge.

10 THE COURT: -- project? Okay. All right. That's
11 helpful. Thank you.

12 From Google, anything to add to the update from Mr.
13 Mullens?

14 MR. SAMPSON: Yes.

15 Good morning, your Honor.

16 THE COURT: Good morning.

17 MR. SAMPSON: Paul Sampson on behalf of
18 Defendants. And just a couple of things that I wanted to
19 kind of run through in response to Mr. Mullens'
20 presentation.

21 In their brief on page five, they cite the Tremblay v.
22 OpenAI case as an exemplar of the number of custodians. In
23 that case, it was 24. We have come forward with the number
24 20 that we think is a good number, a ceiling on sort of what
25 is needed for purposes of where we are in the case. And

1 what I mean by that is, you know, as I think the Court is
2 aware, Judge Lee ordered the parties to focus on class
3 certification issues. And it's been a real struggle to get
4 there. We've been -- you know, I'm sure you saw in a number
5 of the filings that we've had a lot of meet and confers to
6 really try and focus the parties on, "Okay. You've asked
7 for the universe. Let's focus this in on what's really
8 needed for class certification." And the custodians issue
9 is no different, you know, they came in asking for 62, you
10 know, the triple of what's seen in other cases. And those
11 cases are situated differently than us. They don't have
12 class certification coming first. And so that was a big
13 discussion and a lot of meet and confers to really narrow us
14 down.

15 We have identified -- as Mr. Mullens said, we came
16 forward with a compromise. We heard you that we should
17 substantially compromise and to come forward with 20
18 custodians total. So that's seven additional from their
19 list, and we offered, you know, "We'll pick four, you pick
20 three." You know, we're willing to figure out what's the
21 best way to identify that. I don't think it makes sense for
22 them to pick all of them.

23 THE COURT: Why not?

24 MR. SAMPSON: Well, because the -- or the -- the
25 case law talks about the responding party has, you know, the

1 ability to investigate and, you know, determine who the
2 right witness is, and so we've come forward and put forward
3 custodians based on that.

4 In response to the question about Google DeepMind, you
5 know, originally, there were a couple different units within
6 Google, Google Brain, Google DeepMind, that merged, and
7 there's been a lot of reorganization since, and so not all
8 of the products fall within Google DeepMind --

9 THE COURT: Okay.

10 MR. SAMPSON: -- in response to that question.

11 THE COURT: Do the four products that were part of
12 the TRO, were those in DeepMind? It's not a -- it's not a
13 pop quiz.

14 MR. SAMPSON: Yeah, yeah, I don't know the answer,
15 so I'm going to defer to my colleague.

16 THE COURT: All right.

17 MR. TUTTLE: Your Honor, this is Eric Tuttle from
18 Wilson Sonsini.

19 THE COURT: Thank you.

20 MR. TUTTLE: I want to be careful because I don't
21 want to get anything wrong, but I think products -- like,
22 products that are made available to users and public can
23 tend to exist in different organizational units. Google
24 DeepMind, in my understanding, tends to contain the sort of
25 core model development functionalities.

1 THE COURT: Okay.

2 MR. TUTTLE: So there are models that are
3 developed, and then those models are incorporated into --

4 THE COURT: Into products.

5 MR. TUTTLE: -- into products and services. Those
6 products and services, I'm sorry, exist in other
7 departments, organizations, and teams.

8 THE COURT: Okay.

9 MR. TUTTLE: Does that make sense?

10 THE COURT: Yes. Yes, it does.

11 MR. TUTTLE: Okay. Sorry.

12 THE COURT: So I think that answers my question.
13 Thank you.

14 Mr. Sampson.

15 MR. SAMPSON: Yeah.

16 THE COURT: Anything further?

17 MR. SAMPSON: Yeah, another point -- you had
18 questions about search terms and what's being done already.
19 For the 13 custodians, we've already produced a number of
20 documents, applying the search terms. We're, you know,
21 cranking through those. We have 50 document reviewers that
22 we've trained and who are ploughing through those documents
23 at the moment, and the goal is to get documents produced for
24 those individuals before the August 8th class certification
25 motion.

1 THE COURT: Well, it's going to have to be much
2 faster than that, but I understand, and I appreciate the
3 parties' efforts, the discussions, and your efforts even up
4 to and including this morning.

5 Okay. I had reviewed with some care the custodian
6 issues, including the updated joint chart that had both
7 sides' descriptions. I did take judicial notice of and look
8 at the org charts, and I'm -- again, I'm pleased and
9 appreciative of the effort that the parties have made.
10 Where I came out was, again, looking at the products that
11 are likely in the case, and I know you're waiting on the
12 motion to dismiss order, and if there's a -- you know, a --
13 sort of say, "Okay. Well, let's call it six products."
14 There are four that are -- were in from the TRO. There are
15 others that are obviously being debated through the motion
16 to dismiss, if a couple of those get through -- that's how I
17 get to sort of six products. You need a couple of
18 custodians as to each, again in addition to the work at
19 least that had been done at the time I was looking at this.

20 So where I came out in terms of a number, that's my
21 process behind it. It's not completely arbitrary. It will
22 feel that way to each side, I am sure, but the number I came
23 out with was 25 custodians, 25 from the list. If Plaintiffs
24 now want the two declarants, those will count as part of the
25 25. And that gives you -- we've got the 13 now, so that's

1 12 to go, and I think it does make sense for Plaintiffs to
2 have the laboring oar in selection because they have to make
3 the most whittling down, if you will, from their initial
4 ask. I will say, in reading those details, there was
5 certainly a lot of overlap. And I appreciate that the
6 descriptions are necessarily at a fairly high level, but
7 there was a fair amount of overlap, and it certainly
8 appeared to be. And as to certain products, it seemed there
9 was almost a drill down of everyone who had touched it. So
10 I am confident that there is room to further narrow and
11 focus the search, particularly for purposes of class
12 certification. So it will be 25 custodians. That's 12
13 additional, up to the Plaintiffs if they want the two
14 declarants included in that. And beyond that, Plaintiffs
15 can identify a remaining or an additional seven, and Google
16 can identify an additional three. If the Defendant --
17 excuse me, if the Plaintiffs do not want the custodians,
18 then they can use those two picks to further out for the
19 additional -- for the additional 12. All right. So that
20 gets us to a total of 25. Now let's talk about production.

21 MR. TUTTLE: May I ask one question, your Honor?

22 THE COURT: Yes, you may.

23 MR. TUTTLE: Eric --

24 THE COURT: You have to identify yourself for the
25 record.

1 MR. TUTTLE: Sorry, I --

2 THE COURT: That's fine.

3 MR. TUTTLE: Eric Tuttle again from Wilson
4 Sonsini.

5 THE COURT: Thank you, Mr. Tuttle.

6 MR. TUTTLE: I wanted to ask if it would be
7 possible to, within the -- as I understand, we're adding 12,
8 right?

9 THE COURT: Yes.

10 MR. TUTTLE: Within those 12, to have some sort of
11 prioritization by about half for each because we want to
12 make sure to produce the ones that are most critical for
13 class certification before the August 8th deadline. And I'm
14 -- I fear, right, that -- you know, we're well underway on
15 the 13 custodians that have been agreed, but I fear that
16 adding 12 may make it difficult, and I think not all of
17 these custodians I would expect are essential for class
18 certification. Some of them may be more important for the
19 merits, and the discovery period will continue thereafter.

20 So I -- if we could get a -- if we could get an
21 agreement that there was some rank order in terms of
22 prioritization for class certification, that would help us
23 ensure that we can focus discovery efforts on those that are
24 actually needed for -- most urgently in the case.

25 THE COURT: Well, it is in Plaintiffs' interest to

1 prioritize their picks for class cert as well.

2 But let's talk about production --

3 MR. TUTTLE: Okay.

4 THE COURT: -- because it has to speed up
5 substantially, in light of the class schedule or, yes, the
6 class certification schedule that both sides have been aware
7 of and have heard Judge Lee's views on. And we've been
8 through this before with these parties -- excuse me, not
9 with these Plaintiffs, but with Plaintiffs' Counsel with
10 Google, and it is always a challenge, and I appreciate that,
11 because there's a lot of material to review, but production
12 is going to have to speed up, and that may mean you don't
13 have time to review everything that comes out from the
14 search terms, and the parties will have to formulate their
15 stipulation with regards to rights to clawback and
16 qualification under Rule 502(b), but -- and, obviously, you
17 know, I mean that covers any privileged materials, anything
18 that is inadvertently produced, and it will be inadvertent
19 for purposes of 502(b), as well as meeting the other
20 requirements. The reason for that is that, as a practical
21 matter, it is -- there simply may not be time to review as
22 the documents are produced. And I don't know what the
23 custodians have and what the volume is, and I know that that
24 makes Defendants uncomfortable, but it is the practical way
25 that production has to proceed. And with the appropriate

1 stipulation and safeguards in place, then that should --
2 that shouldn't be a problem. If it's also as it relates to
3 technical documents -- and, again, if you need to produce
4 prior -- without review and you set that at your highest
5 level of priority, subject to anything can be de-designated
6 or lower designated later, then that is -- that's fine as
7 well, because the Defendant needs these protections in
8 place. But this production is going to have to get
9 underway. I'm glad that it's already -- sounds like
10 substantially underway as to the 13, so as to these
11 additional, it's -- we've got to get it going. So rolling
12 production will have to get started. We are already there,
13 so -- we have to get started in two weeks and continue every
14 two weeks. And the due date on the class cert motion is the
15 11th?

16 MS. WEAVER: I believe that's right, your Honor.

17 UNIDENTIFIED SPEAKER: I believe it's --

18 MS. WEAVER: I'm sorry, is --

19 UNIDENTIFIED SPEAKER: -- August 8th.

20 THE COURT: August 8th --

21 MS. WEAVER: August 8th, yes. Yeah, and we do
22 have a 502(d) stipulation already in place, your Honor, and
23 I would --

24 THE COURT: Okay.

25 UNIDENTIFIED SPEAKER: (Indiscernible).

1 THE COURT: This is Lesley Weaver for the
2 Plaintiffs.

3 MS. WEAVER: I apologize.

4 THE COURT: Okay.

5 MS. WEAVER: Okay. Yes, we will work on a 502(d).
6 I apologize for the misstatement.

7 THE COURT: That's fine.

8 MS. WEAVER: To date, we have received 5,000
9 documents --

10 THE COURT: Okay.

11 MS. WEAVER: -- and we produced 10,000 for the
12 Plaintiffs, so we're going to move. And if your Honor -- I
13 think, just to clarify, are you setting date production
14 deadlines?

15 THE COURT: Yes, rolling production. We'll get
16 started within two weeks and continue with the production
17 every two weeks until the 29th of July, it'll need to be
18 done, as to class cert. And you do owe the -- I'm speaking
19 to the Plaintiffs. You do owe the Defendants a
20 prioritization on these -- on the custodians.

21 MS. WEAVER: And we can do that. Is there a
22 deadline by which you would like us give them the
23 prioritization and --

24 THE COURT: Well, it certainly behooves you to do
25 it as soon as possible.

1 MS. WEAVER: We'll do it this week. Okay. Thank
2 you.

3 THE COURT: All right.

4 MR. TUTTLE: May I just be heard on two things
5 related to that, your Honor.

6 THE COURT: Yes, Mr. Tuttle

7 MR. TUTTLE: So I understand the direction --

8 THE COURT: That's right. Mr. Tuttle.

9 MR. TUTTLE: Mr. Tuttle speaking, sorry.

10 THE COURT: That's okay.

11 MR. TUTTLE: I can't remember to do it.

12 THE COURT: Well, it's hard to have to keep saying
13 it, but I appreciate the effort. Thank you.

14 MR. TUTTLE: I guess it means I'm getting up too
15 often, your Honor, and I apologize for that.

16 But I just want to be heard on two things. One is we
17 certainly hear the direction to get the -- we have been
18 producing. We will make sure we continue rolling every two
19 weeks and to get done by July 29th. I think -- I've got a
20 lot of concern, I think naturally, with producing documents
21 without a privilege review, and I think there would be
22 issues with an order requiring us to do that -- I -- so I
23 think we will look at the best way to accomplish the
24 schedule that the Court has directed us to do. I do not
25 understand the Court to be ordering us to produce without a

1 review. I understand the Court to be ordering us to figure
2 out how to do this.

3 THE COURT: How to do it on this timeframe. And,
4 often, what happens in this case is, as part of the review,
5 you search for, obviously, counsel name in the documents or
6 reference, but page by page review just won't be feasible on
7 this schedule. And, again, the protections are in place, or
8 the parties will get those in place, with regards to the
9 502(b) issues, so do your best.

10 MR. TUTTLE: We will figure out a way to get it
11 done.

12 The other thing is, I just want to -- I don't want
13 there to be a -- this needs to go both directions, I assume.
14 And Ms. Weaver talked about the total number of documents
15 that have been produced, but we've run through the
16 Plaintiffs' -- you know, Steve Almond has produced 32
17 documents, which are 12 e-mail, so --

18 THE COURT: Those issues aren't in front of me
19 today.

20 MR. TUTTLE: I understand. So we just -- but to
21 be clear, this is -- everybody needs to get this done.

22 THE COURT: Everybody needs to get it done if it
23 relates to class certification, which I imagine you have
24 made some requests that go to that, directed to the
25 Plaintiffs.

1 MR. TUTTLE: Yes.

2 THE COURT: And those documents need to be turned
3 over.

4 MR. TUTTLE: Okay. Thank you, your Honor.

5 THE COURT: You're very welcome.

6 MS. WEAVER: And I apologize, but to clarify for
7 the record, the protective order refers to the clawback
8 procedure, and it was entered on April 4th. That's docket
9 number 119 and paragraph 15 is -- deals with clawback and
10 Rule 502(d). But we don't have a 502(d). We can endeavor
11 to get that in place as soon as possible.

12 THE COURT: We have excellent lawyers on both
13 sides. I'm confident you can meet and confer and get that
14 done.

15 MS. WEAVER: Great.

16 THE COURT: Okay. That takes care of custodians.
17 Now let's turn to the easy issue, shall we? So let me share
18 with you my initial thoughts on source code. I don't have
19 quite enough room up here, so let me get everybody
20 organized. Okay. Source code.

21 First, let's address the rightness of the source code
22 request. I reviewed the request for productions, the
23 definitions that were used, the arguments on both sides as
24 to what was included in which requests. The initial RFPs do
25 not ask for source code where examples of documents that are

1 included in the requests are given. And there are many
2 examples of documents and communications, but source code is
3 never listed or addressed there. Obviously, the second set
4 of RFPs does expressly address source code. It was designed
5 for that. Google has made its objections to production of
6 source code very clear. Our record is robust. And so for
7 judicial efficiency, we're going to proceed and address it
8 now. I'm not quite sure where the response date was, if
9 we're still ahead of it.

10 UNIDENTIFIED SPEAKER: We're still ahead of it.

11 THE COURT: That's what I -- that's what I
12 thought. But we are going to go ahead and address the issue
13 now. And I'm taking Google stated objections with regards
14 to this dispute in the joint statement and the declarations,
15 et cetera, as -- that's the record, and I'm taking those as
16 objections. They are noted, and we're going to work our way
17 through those here today.

18 Where I always start in these cases is with the class
19 definition, because, of course, when we're talking about
20 source code, I am looking at relevance, and I'm looking at
21 necessity. As I understand it from Judge Lee's order where
22 she was striking the fail safe allegations, she said, you
23 know, here's a definition that works, "All persons in the US
24 who own a US copyright that was used by Google to train
25 Google's generative AI models during the class period." If

1 there's something else that the parties are working from, I
2 want to hear about that. And Mr. Tuttle is shaking his head
3 no on behalf of Google.

4 And, Plaintiffs, is that our working class definition?

5 MR. YOUNG: I believe so, your Honor, with the
6 exception --

7 THE COURT: Identify yourself, please.

8 MR. YOUNG: Oh, sorry. Christopher Young, Joseph
9 Saveri Law Firm, for Plaintiffs.

10 That sounds about right, with the caveat that we've
11 endeavored to limit the discovery that we're seeking,
12 including the -- the source code we're seeking to inspect,
13 to certain models and exemplar versions of those models.

14 THE COURT: I understand, but this is our working
15 class definition.

16 MR. YOUNG: Correct.

17 THE COURT: Is that correct?

18 MR. YOUNG: Correct.

19 THE COURT: Okay. And what is the class period?

20 MR. YOUNG: We contend that it begins from 2017.
21 The precise month and date, I don't recall offhand, but
22 there is, I think, a putative or pending dispute about that
23 temporal limitation. I understand that my friends at Google
24 might contest whether or not that is an appropriate class
25 period, but we -- as alleged, it's 2017 onwards.

1 THE COURT: Okay. All right. We'll start with as
2 alleged. All right.

3 And the products that I understood, again from Judge
4 Lee's tentative, where she was was to allow for Bard,
5 Gemini, Imagen, and -- how do you pronounce LaMDA
6 (pronouncing)?

7 MR. YOUNG: LaMDA (pronouncing).

8 THE COURT: LaMDA (pronouncing), thank you. Like
9 the lambda.

10 MR. YOUNG: Like the Greek letter.

11 THE COURT: Like the lambda. Okay. Got it.
12 Thank you.

13 And then others were somewhat on the fence. And I
14 appreciate or I noted the Plaintiffs' request from those for
15 all the various versions, which seems to expand the list
16 substantially. I understand the thinking behind that.

17 So where I start is, it certainly appears that a
18 relevant question for class certification is, what was
19 Google sourcing -- harvesting from the internet? And we're
20 obviously focused on data from the internet in the relevant
21 time period. And how -- and then, ultimately, how was that
22 used?

23 So my -- you know, of course, the overarching question
24 is, is source code the best indicator of that? And I've
25 read through the declarations taking issue with where URLs

1 are shown, et cetera, but my first question is, really, what
2 has Google produced to date that shows the -- well, we'll
3 start with shows the content of training sets. Now, I know
4 we have the -- what was harvested from the internet, and
5 then what actually ends up in training sets, and then how
6 those are used. But what has Google produced to date that
7 shows what actually -- what -- the content, if anything,
8 that shows content of training datasets?

9 And, Mr. Wong (sic), respectfully, let me get Google's
10 answer. And I'm going to hear argument from both sides, but
11 I want to get some basic elements in before me.

12 MR. YOUNG: Yes. And, your Honor, just to nip
13 this early before it -- it's Mr. Young. Thank you.

14 THE COURT: Well, Mr. Young, I apologize. I
15 apologize.

16 MR. YOUNG: No problem.

17 THE COURT: I'm going to put my Plaintiffs where
18 my Plaintiffs belong. Okay.

19 MR. TUTTLE: Thank you, your Honor.

20 THE COURT: Mr. Tuttle.

21 MR. TUTTLE: So the --

22 THE COURT: Mr.?

23 MR. TUTTLE: Sorry. I will get this. Mr. Tuttle
24 speaking.

25 The -- as we've -- we've had a number of -- we've had

1 many discussions with the Plaintiffs about the training
2 data, and basically it's an extraordinarily large amount of
3 data. And so we proposed from the outset that what we do is
4 we provide the Plaintiffs with information describing the
5 datasets used for each model, and then Plaintiffs select
6 representative or sample datasets. We eventually agreed at
7 some level on doing so, and we got a -- we got a set from
8 them. We provided these data cards that -- for models, for
9 all the different models and all the different --

10 THE COURT: So what's a data card? What's on the
11 data card? We're drilling down from datasets, I take it.

12 MR. TUTTLE: Well, I'm going to get to the
13 datasets, but the --

14 THE COURT: Okay.

15 MR. TUTTLE: -- data card, you know, I'm using
16 that as a general term. They vary, because we're talking
17 about a long period of time and different teams, but at a
18 high level, a data card describes a model and the training
19 datasets that were used for that model.

20 THE COURT: Okay.

21 MR. TUTTLE: So it describes in English, right,
22 with descriptive names. And the amount of detail and
23 additional information varies, right, by model and time
24 period and team.

25 THE COURT: So you said it identifies a model and

1 the training data used for that model?

2 MR. TUTTLE: Correct. It identifies the -- just -
3 - to be clear, it identifies sets of data, so, you know, web
4 -- like, it might identify a set as Wikipedia, or it might
5 identify a set as data crawled from the internet, or it
6 might -- you know, and so on. And then it will -- it will -
7 - and then it will -- and it'll have additional information.
8 So we provided this to the Plaintiffs, so that they could
9 select datasets that would be of interest to them, because
10 there are, we assume, some datasets that are of more
11 interest to the Plaintiff than others. The Plaintiffs made
12 a selection, and I don't have a date offhand --

13 UNIDENTIFIED SPEAKER: May 20. May 20.

14 MR. TUTTLE: May 20th, we were given the datasets
15 that Plaintiffs selected.

16 THE COURT: Okay.

17 MR. TUTTLE: And we have -- we -- this is another
18 -- whole other issue that I don't want to spend time on, but
19 we -- it's too large, it's extremely large, and it's also
20 very sensitive. So what Google offered to do is to make it
21 available in a review environment within Google --

22 THE COURT: A clean room environment --

23 MR. TUTTLE: Well, it's a remotely accessible
24 environment. It uses Google's cloud technology, right, to
25 make an environment in which this data will be placed and

1 where the Plaintiffs and their experts, under the protective
2 order, can securely connect remotely using a Chromebook.
3 And they will be able to access this environment where all
4 the data will be, and they will be able to operate on it and
5 explore it, right? And --

6 THE COURT: And you're talking about just one of
7 the selected datasets?

8 MR. TUTTLE: All of the selected datasets --

9 THE COURT: Okay.

10 MR. TUTTLE: -- will be -- that can -- you know,
11 that -- all of the selected datasets that, in fact,
12 correspond to a model -- like, we're working through it. I
13 think in some cases, they may have picked a dataset that was
14 not actually used with any of these models. But to the
15 extent the dataset was correctly identified as a dataset
16 used with the model, the dataset is being located, is being
17 moved into this environment, so that the Plaintiffs can
18 review it.

19 THE COURT: Uh-huh.

20 MR. TUTTLE: Even this limited number of datasets
21 is very, very large, so we are in the process -- like, I
22 think about half of it is currently in there, and the
23 environment has been built. This was slowed down a bit by
24 disputes between the parties about the nature of the
25 environment, but we went ahead -- given the timing, we just

1 went ahead and built it and put the data in there. And I
2 think we have -- we'll be shipping the Chromebooks to the
3 Plaintiffs this week, and we just need to get the account
4 information, and I think we should be able to get them
5 access, so --

6 THE COURT: And they can access it remotely?

7 MR. TUTTLE: They can access it remotely. They
8 can access it from anywhere. It will be -- I mean, it will
9 be -- they need to be on one of these laptops.

10 THE COURT: I understand.

11 MR. TUTTLE: They need to be one of these
12 credentialed users. And there are various security measures
13 in place. So the data is being made available there. And I
14 think the Court is quite right to focus on the training
15 data, because the training data is going to tell us what
16 works are -- were used to train the models.

17 THE COURT: So -- and will the training data, does
18 it identify specific works or specific content?

19 MR. TUTTLE: The training data contains the
20 content.

21 THE COURT: Okay.

22 MR. TUTTLE: And, typically, it contains metadata
23 -- and, generally, it contains metadata about the content,
24 which -- for example, for web crawl data that Google
25 crawled, it should generally contain the URL. There are,

1 you know, thousands and thousands and thousands of datasets
2 across many teams over many years, so we haven't been able
3 to look at all of it and determine if that is always the
4 case. And if it should turn out that there's information
5 that is not present, we can figure out a way to deal with
6 that. But that information, in any event, will not be in
7 the source code.

8 THE COURT: But your understanding, your
9 representation to the Court is that the datasets reflect the
10 training model content --

11 MR. TUTTLE: -- the content.

12 THE COURT: -- the content and the metadata.

13 MR. TUTTLE: Generally.

14 THE COURT: Generally. So why generally? Or just
15 tell me what your --

16 MR. TUTTLE: I say that only because -- like,
17 because we're -- it clearly contains the content, right?

18 THE COURT: I got that. I got that.

19 MR. TUTTLE: Because that's the content that is
20 exposed --

21 THE COURT: I'm trying to get to the URL. So tell
22 me how --

23 MR. TUTTLE: Yeah. So the URLs for, again -- what
24 -- my understanding, we just haven't been able to review all
25 of the training data, so -- because it's many, many -- as I

1 said, it's many, many thousands of files all over the place,
2 but my understanding is that for -- as for, like, public web
3 stuff that Google crawled and collected and used for
4 training data, that the URL metadata should generally exist
5 in -- the URL metadata should generally exist in the
6 training data files. To be clear, these are the files that
7 are the -- that are -- the training data that is available
8 to the model to use to train.

9 THE COURT: Uh-huh.

10 MR. TUTTLE: There is some sampling that occurs
11 during the training and so forth. And so it's -- it -- we
12 can't say -- like, there are -- in some cases, we -- like,
13 the -- what was actually used got saved, and so we would
14 have a record of that. But in other cases, all we have is
15 the final set of the data that was available to be trained
16 on, if that makes sense.

17 THE COURT: What's the difference between those
18 two?

19 MR. TUTTLE: So it could be that the training
20 datasets that the model draws from to train --

21 THE COURT: Yes.

22 MR. TUTTLE: -- there could be data in there that
23 was not sampled, but the data that could have been chosen is
24 in those files.

25 THE COURT: Okay. All right. And you say you are

1 shipping Chromebooks this week?

2 MR. TUTTLE: Yes, I think they can go out.

3 UNIDENTIFIED SPEAKER: We just need to know who to
4 ship --

5 MR. TUTTLE: Yeah, we just need to get the
6 information on --

7 THE COURT: All right.

8 MR. TUTTLE: -- who to ship them to.

9 THE COURT: So -- and is there any reason why
10 review in this environment, or, excuse me, why Plaintiffs'
11 access to this environment can't be up and running next
12 week?

13 MR. TUTTLE: I'm not aware of a reason. Oh, well
14 -- there is an ongoing dispute about cost, and I --

15 THE COURT: Okay.

16 MR. TUTTLE: -- it has been teed up for a motion,
17 but I don't know that it is necessary to resolve because, in
18 our view, it is premature.

19 THE COURT: Okay.

20 MR. TUTTLE: Google has agreed to set this up at
21 its own expense.

22 THE COURT: Okay.

23 MR. TUTTLE: And -- but it has -- and it has said
24 that -- you know this environment --

25 THE COURT: I'll look forward to the motion --

1 MR. TUTTLE: Yes.

2 THE COURT: -- if the parties are unable to work
3 it out.

4 MR. TUTTLE: Yes.

5 THE COURT: But it won't stop the Plaintiffs from
6 accessing as early as next week?

7 MR. TUTTLE: No, it will -- from our perspective,
8 it does not. We have offered to make it available up to --
9 I believe it's \$1 million of commercial rates of free use.
10 We did reserve the right to pause it if they go well past \$1
11 million, but unless they hit \$1 million in commercial rates,
12 there's nothing -- there won't be anything stopping them
13 from using it, from our perspective.

14 THE COURT: Okay. Well, as I say, I don't expect
15 access to be limited on any cost issue at this time, and,
16 obviously, if the parties can't come to an agreement, you'll
17 come back to me, and we'll deal with it.

18 MR. TUTTLE: Yes. It's premature, because I hope
19 we will never get to \$1 million.

20 THE COURT: I hear you. I hear you. I hear you.
21 I hope that's right, but we're talking about a lot of data,
22 so different issue.

23 Okay. So I asked you to come to the podium, Mr.
24 Tuttle, on behalf of Google, because my question was, so
25 what has been produced? And it sounds like that -- again,

1 I'm pleased the parties had discussions and went to the
2 sampling approach and that Plaintiffs identified datasets,
3 and Google has endeavored to make those available through
4 this environment, and Plaintiffs' access and inspection of
5 that should begin as early as next week.

6 MR. TUTTLE: Okay.

7 THE COURT: Okay. Don't go away. Let me just --

8 MR. TUTTLE: I'm right here, your Honor.

9 THE COURT: Okay. That's very helpful. Thank
10 you, Mr. Tuttle.

11 MR. TUTTLE: Thank you, your Honor.

12 THE COURT: Mr. Young.

13 MR. YOUNG: Thank you, your Honor.

14 THE COURT: Thank you for making sure I'm calling
15 up the right person. It's your motion. Let me start -- let
16 me hear from you. I appreciate the explanation from Google.
17 That seems to get us a significant way down along the lines
18 of your request for production 65. So, again, with the
19 necessity for source code standard firmly in mind, what's
20 the issue?

21 MR. YOUNG: Thank you, your Honor. So before I
22 talk about why data cards --

23 THE COURT: Pull a microphone close to you, Mr.
24 Young, just so we're sure we're getting the recording.

25 MR. YOUNG: Thank you. Is this better?

1 THE COURT: That sounds good.

2 MR. YOUNG: Great. So before I talk about the
3 source code, I do want to respond to a couple of points
4 about the training data and the training data protocol,
5 because I would be remiss if we just didn't put that in the
6 record. I mean, from our position --

7 THE COURT: Just slow down a little bit.

8 MR. YOUNG: Yeah. From Plaintiffs' position, one
9 of the reasons why this has being, like, teed up for a brief
10 at all is we've taken issue with how the training data has -
11 - is going to be produced to us, just given that we've
12 engaged in months of -- months of negotiations, and we feel
13 like Google has foisted unilaterally their initial position,
14 essentially wiping away months of effort, including on this
15 cost issue. But I understand that that issue is not ripe,
16 but I just would be remiss if we just let those statements
17 stand without some sort of response here on the record.

18 THE COURT: Okay.

19 MR. YOUNG: Now, with respect to the question
20 about why data cards and model cards are not enough, what
21 that really answers is the what, not the how and the why.
22 With respect to categories of source code we think are, you
23 know, I wouldn't say essential, but are the best evidence of
24 some of the conduct at issue, that will need to be shown at
25 a class wide -- on a common class wide basis, or the source

1 code, as your Honor identified, have the organization and
2 collection of data, so crawlers scraping from the internet,
3 the processing of that data, so processing from raw data
4 into a form that is ingestible by the machines, by the
5 models for the training process, and at the end step, what
6 is needed to prevent what is called leakage or regurgitation
7 of training data from the model. And, you know --

8 THE COURT: So let's focus on what source code --
9 what your argument is as to the need for source code for
10 purposes of class certification.

11 MR. YOUNG: Right.

12 THE COURT: Let's put aside the issue, to the
13 extent we can, of fair use, and, you know, what I think of
14 that is what happens to data once it's obtained. But if you
15 have access to, "Here's the training data" and you're able
16 to review that, again keeping the class definition in mind -
17 - you're able to review that, you're looking for text or
18 images or whatever that is subject to copyright and is
19 included in a training set, you're also able, at least
20 generally, using Mr. Tuttle's word -- should be able to look
21 at the metadata as to where that came from. So if you have
22 a URL and you have the data that ended up in the training
23 set whatever -- in whatever form, doesn't that inform your
24 issues regarding class certification for numerosity,
25 typicality, and et cetera?

1 MR. YOUNG: Perhaps. And I say that because the
2 data cards and the model cards we have received are
3 essentially snapshots. The source code would really be the
4 best evidence of all that data. And I point your Honor to
5 paragraph 11 of the Liekti (phonetic) declaration that was
6 submitted by Google. I mean, for example, that paragraph --

7 THE COURT: Hold on. Hold on.

8 MR. YOUNG: Yes.

9 THE COURT: Let me get there, Mr. Young.

10 MR. YOUNG: Sorry, your Honor.

11 THE COURT: That's all right. I'm going to keep
12 slowing you down, so --

13 MR. YOUNG: Yeah. So whenever you're ready, your
14 Honor.

15 THE COURT: Wait, let me -- I read this with some
16 care. Where did it go?

17 MR. YOUNG: So Mr. Liekti -- and I believe this is
18 the second sentence of paragraph 11 of his declaration --
19 says,

20 "The source code for training generally
21 includes mixture files, so scripts that
22 identify the file names and locations of
23 the various compilation datasets that
24 are used in the training process, as
25 well as proportions by which these

1 datasets are sampled from."

2 So why is that important for class certification,
3 right? What Mr. Liekti is saying, if I'm understanding him
4 correctly, is that their source code that compiles the
5 various training datasets that Google maintains into what is
6 the -- what is going to be the training corpora of that
7 model, right? So whether or not a training dataset is used
8 or is not used for a training dataset, this code will tell
9 us what exactly is used. Further, it will tell us what
10 proportion that is being used, but we anticipate that one of
11 the things Google may say is that a dataset may not -- a
12 dataset containing copyrighted information may or may not
13 have been used more or less, right, which may have
14 implications, not only for fair use, but also for class
15 certification and creating -- perhaps create individualized
16 issues for datasets. So that's why we would need
17 information about that source code. Right.

18 Setting that aside --

19 THE COURT: Hang on. Hang on.

20 MR. YOUNG: Yeah.

21 THE COURT: Let me just take a look at the
22 language. Okay.

23 MR. YOUNG: Okay. So, you know, your Honor, if
24 you reflect back on the, kind of, three-part framework of
25 the three categories of source code identified early on in

1 my argument that we would argue are -- we would most like or
2 would be most explanatory for class certification purposes,
3 this is really the organization or processing point, because
4 this is the organization or collating of data before it
5 becomes used as training data for the models, right?

6 Now, taking a step back to the web crawlers. Now why
7 is that important? The web crawlers, including kind of the
8 filtering -- and this is described by Mr. Liekti in
9 paragraph four of his declaration, as well as by Mr. Ibrahim
10 in his declaration, identifying specific paragraphs.

11 THE COURT: I'm sorry, Mr. Young, you need to slow
12 down.

13 MR. YOUNG: I apologize.

14 THE COURT: That's fine.

15 MR. YOUNG: I will take a breath. This is a bit
16 of a problem for me, so I will endeavor to correct my
17 cadence.

18 Now, if you'll --

19 THE COURT: So we were talking about -- you wanted
20 to move back, if you will, to crawler source code.

21 MR. YOUNG: That's right. And the reason I am
22 framing it this way is it's sometimes helpful to think about
23 the life cycle of data as it gets fed into the model, right?
24 So the crawling is really the first step --

25 THE COURT: Right.

1 MR. YOUNG: -- so this is the original copy. So
2 besides identifying class works, what is -- why this is also
3 important for class certification is because it also bears
4 on common issues of intent and willfulness on behalf of
5 Google, because the training dataset will not tell us what
6 processes Google may have done at that initial collection
7 stage to filter, include, or exclude data, right? You could
8 imagine, right? And, you know, of course we haven't seen
9 the source code, but even from a common understanding of how
10 Google works, Google Search, which Mr. Liekti identifies in
11 paragraph four as being the provenance of training data, has
12 -- blocks or restricts access to certain websites. If
13 Google has chosen to restrict or not chosen to restrict, for
14 example, pirate websites that it crawls, that is going to be
15 common evidence used by the class to -- that would be common
16 evidence used by the class that goes directly to -- only to
17 questions of infringement.

18 THE COURT: But if you have the URL data in the
19 metadata from a training set, then aren't you aware of
20 what's coming from a -- you know a legitimate website or a
21 pirated website?

22 MR. YOUNG: Well, we also -- the issue is also we
23 do not know what the -- how Google maintains its training
24 data. For -- from my experience, having litigated other
25 generative AI cases, it can be a mixed bag, but, generally -

1 - but we do know that the source code extracts URLs and text
2 from the -- from those sources that Google crawls. And I
3 would point you to Mr. Zhao's declaration. I think this is
4 paragraph three of his declaration, where he admits that
5 Google's crawlers, one of the functionalities is to extract
6 URLs and content that are destined to be -- for use as
7 training data. So we know that the source code evidence is
8 that, and the -- well --

9 THE COURT: That's hotly disputed as to whether or
10 not you can tell URLs from the source code.

11 MR. YOUNG: Yeah. Well, I would --

12 THE COURT: So tell me where you think that the
13 Zhao declaration points to that?

14 MR. YOUNG: Yes. I will point you to paragraph
15 three --

16 THE COURT: Uh-huh.

17 MR. YOUNG: -- and the second sentence beginning
18 "these crawlers." So --

19 THE COURT: Hang on. Hang on. Hang on. It's
20 actually the third sentence, but that --

21 MR. YOUNG: I apologize. The little -- this
22 little second sentence kind of evaded my eyes there.

23 THE COURT: All right. I see the language about
24 extracting URLs.

25 MR. YOUNG: Correct. You know, there's an

1 inconsistency I would admit within the Zhao declaration
2 wherein he says in paragraph five -- and your -- if your
3 Honor would -- I mean, I would be remiss if I didn't point
4 that out -- say that the source code will not identify
5 specific URLs, but this seems in direct tension with
6 paragraph three of that declaration, where he says that it
7 would extract URLs. And even so, if the -- if -- what
8 Google has asked us to do is to pick datasets essentially in
9 the dark with imperfect information, whereas the source code
10 would give us a direct snapshot into what is used to
11 organize and collate training datasets. We are being asked
12 to select representative datasets from data cards and model
13 cards, which are snapshots of various models and versions
14 which we have very limited information about, whether
15 they're production models, benchmark models, intermediate
16 models. But the source code was really the best source of
17 evidence for us to suss out what relevant datasets we really
18 want and really need, and --

19 THE COURT: How would you possibly have time for
20 that, given the vast amounts of code that you're talking
21 about?

22 MR. YOUNG: So it is -- so, you know, there's a
23 bit of a tension here, right, because our experts do not --
24 like, we -- taking at face value, right, that Mr. Liekti is
25 correct that there are millions of lines of code, my

1 understanding from our experts is that a lot of that code is
2 not any code that we are interested in at all. I mean,
3 really from even just these declarations -- and I think our
4 experts have identified, just from the scant document
5 production, the code -- specific code that we are actually
6 interested in, which I don't think Mr. Liekti takes much --
7 no, he quibbles with it, but I think they're actually more
8 in line than they are in disagreement. We have identified
9 specific code that we are interested in. For example, to
10 the extent that they're -- you know, you would believe --

11 THE COURT: With what specificity when you say
12 you've specified?

13 MR. YOUNG: Yeah. So Mr. Ibrahim has gone through
14 some of the documents that Google has provided, with the
15 understanding that was before the recent, kind of, custodial
16 -- first wave of custodial documents, and identified
17 specific scripts and code that were identified in the
18 documents. I mean, he has identified something he believes
19 is a scrape code, which we understand from Mr. Liekti's
20 declaration to be something that Google uses to convert
21 crawled data into training-ready data, right? And I don't
22 think Mr. Liekti has disagreed with what Mr. Ibrahim says --
23 said what this is at all, but, really, rather than being
24 scraped from the internet, it was scraped from something
25 that was already corrected -- collected from the internet,

1 excuse me. Mr. Ibrahim has also identified code that
2 identifies leakage and regurgitation and how that would
3 "require indexing training data," which is why we're also
4 interested in that last stage, because that indicates
5 Google's knowledge and knowledge of whether or not the
6 training datasets that it collected contain copyrighted
7 works, which would again be common evidence that each member
8 of the class would rely on. But without actually knowing
9 what's in the code and what's referenced in the code, even
10 if it's not looking at referencing specific works in the
11 training dataset, just knowing, for example, "This is the
12 dataset that we are indexing, and that's the dataset we want
13 to look at," right, setting aside, you know, there might be
14 a dataset, for example, called, like, e-mails, which we
15 might not be interested, right? But just having the
16 references or the identification of the datasets that Google
17 is using in its source code and referencing in its source
18 code would be tremendously helpful at class certification.

19 Now, your Honor was giving me a bemused look, so I will
20 pause to see if I can offer some clarity.

21 THE COURT: No, I'm following your argument. I
22 get the -- I get that you want the source code. I
23 understand. I hear the argument of, "Well, it's the best
24 evidence of." I'm trying to get to the necessity point.
25 And, again, I'm focused on the class certification issues

1 and the identification of where did Google go to get data to
2 build training sets? And it sounds -- and, again, if the
3 datasets in the -- that are available for inspection contain
4 the metadata, it seems like that goes a long way to
5 identifying where they went, what they collected, and then,
6 you know, obviously, there's issues around what actually
7 ends up in the training data.

8 MR. YOUNG: Yeah. Well, your Honor, I think just
9 to directly respond to that point, right, I mean, we have --
10 we have been -- you know, as -- I believe I heard Mr. Tuttle
11 say, right, that even one of the representative datasets
12 that we've identified from the data cards and model cards
13 may not have been used by Google at training, which really
14 highlights why the process we've been doing here now is an
15 imperfect one, because we are being asked again just to kind
16 of pick and choose, from the limited information we have,
17 what the most relevant training datasets would be for class
18 certification. And part of the difficulty we've been having
19 is that we do not know for a fact which -- or know with any
20 reasonable amount of certainty which datasets have been used
21 for training the models or if any of the datasets themselves
22 contain copyrighted works at all, because they are --
23 generally, they're not descriptive. We do not know like --
24 besides, like, the name --

25 THE COURT: No, but you're going to have a chance

1 to investigate --

2 MR. YOUNG: Of course. Of course.

3 THE COURT: -- them and look behind the curtain,
4 if you will. Have the Plaintiffs asked or endeavored to ask
5 Google as to when you got to the -- went through the
6 sampling and selection of datasets, as -- was there a
7 discussion with Google for confirmation as to whether or not
8 the datasets identified by Plaintiffs were actually used in
9 training?

10 MR. YOUNG: There's been a long back and forth
11 about the training datasets. Even getting to the point of
12 getting data cards in the first instance was a suggestion
13 from Plaintiffs. So, your Honor, we've been working really,
14 really hard to attempt to identify the training datasets.
15 And, you know, this is the first for -- at least from my
16 recollection, and perhaps my colleagues have heard this
17 before, that even one of our own training datasets was not
18 used. We've been simply asked to identify training
19 datasets, and we provided them a list.

20 THE COURT: I'm sorry?

21 MR. YOUNG: We've been asked to provide a
22 representative set of -- representative datasets --

23 THE COURT: Right. And have you asked Google to
24 confirm that those datasets were actually used in training?

25 MR. YOUNG: Not -- no, I do not think we have --

1 THE COURT: Okay.

2 MR. YOUNG: -- but we -- based on -- I mean, based
3 on the model dataset -- the data cards that we've had, these
4 are the most -- at least, based on the names of the
5 datasets, these seem to be the most likely ones to be --
6 contain class works.

7 THE COURT: All right. That's very helpful, Mr.
8 Young. Thank you.

9 MR. YOUNG: Uh-huh. Now, your Honor, I do want to
10 -- we've talked about crawlers and the, kind of,
11 organization filtering scripts. I do want to touch briefly
12 on the last category of scripts that I was referencing,
13 which is --

14 THE COURT: Okay. Hang on. Hang on. Hang on.

15 MR. YOUNG: Yes.

16 THE COURT: Give me just a second. Okay. Go
17 ahead.

18 MR. YOUNG: Yeah. So the last category of source
19 code that we are seeking are these "mitigation-type
20 filters."

21 THE COURT: Uh-huh.

22 MR. YOUNG: So these are scripts at the back end
23 that check, essentially, for lack of a better word, what is
24 coming out of the model. So, for example, this could be
25 checking for PII, for example, if someone's, like, e-mail

1 address or someone's, like, address is about to be spit out
2 by the model, it'll check. Why we think this is relevant
3 for class certification is that this is just another source
4 of relevant training data, because, for example, if --
5 because how we understand these things to work, based on the
6 limited discovery we have, as well as just our experts'
7 understanding -- broad understanding of how these generative
8 AI models function, you cannot have one of these filters on
9 the back end unless it references what you are trying to
10 stop from coming out. So you can
11 imagine --

12 THE COURT: Well, when you say on the back end --

13 MR. YOUNG: Yes.

14 THE COURT: -- where is it? Is it on the training
15 data? Is it as the data moves from the crawler into the
16 training dataset? Where?

17 MR. YOUNG: So this is on the very -- right before
18 the model outputs information --

19 THE COURT: Uh-huh.

20 MR. YOUNG: -- for example, you could imagine the
21 model just taking out of the copyright -- or maybe keep it
22 in, right? It is about to emit, for example, two, like,
23 sentences for Harry Potter or a chapter of Harry Potter.
24 The filter will go, "Hey, wait, this is a copyrighted book.
25 Don't do that. Instead, say this," right?

1 THE COURT: Uh-huh.

2 MR. YOUNG: Say, "I cannot emit copyrighted
3 works."

4 THE COURT: Okay.

5 MR. YOUNG: Now, how we understand this to
6 generally work -- and, you know, perhaps Google has devised
7 a better way to do it, but based on the document that Mr.
8 Ibrahim has identified that says that Google needs to index
9 training data to develop this, we do not think that they do.
10 That means that they have to develop the filter using copies
11 of the copyrighted work that they are trying to stop from
12 coming out --

13 THE COURT: Okay.

14 MR. YOUNG: -- which, of course, is common
15 evidence of knowledge of -- that their original training
16 corpus contains copyrighted work.

17 THE COURT: So how does that relate to the need
18 for source code?

19 MR. YOUNG: The -- because this is evidence in the
20 source code. It'll refer back to the -- either the training
21 dataset or some training dataset of a collection of
22 copyrighted works. And just the mere existence of this is
23 also evidence of that intent. But the source code would be
24 the best evidence of the specific concerns and intent that
25 Google had about stopping copyrighted works from being

1 emitted.

2 THE COURT: Let me make a note here. Let me look
3 at something. Thank you, Mr. Young. Let me see if I have
4 another question for --

5 MR. YOUNG: Absolutely, your Honor. Thank you.

6 THE COURT: Okay. All right. That's very
7 helpful. Thank you.

8 MR. YOUNG: Thank you very much, your Honor.

9 THE COURT: All right. From Google, I have some
10 follow up questions. Who's got the mic on this one?

11 MR. TUTTLE: Thank you, your Honor. Mr. Tuttle.

12 THE COURT: Mr. Tuttle for Google.

13 So let's go to a couple of my questions. I will give
14 you an opportunity to kind of backtrack and address the
15 explanations provided by Mr. Young. So datasets -- sample
16 datasets identified by the Plaintiffs. You've been working
17 with those to get them in this environment. Has there been
18 confirmation by Google, or is Google able to confirm whether
19 or not those datasets were actually used in training?

20 MR. TUTTLE: Yes, we are able to do that, your
21 Honor, and we do intend to -- we're working through the
22 list, and we're intending to, you know, explain to the
23 Plaintiffs what we've found. I think where the Plaintiffs
24 used a data card for a specific model and chose a dataset
25 described on that data card, that dataset was used, and we

1 have it. I think there were a few instances where
2 Plaintiffs chose to use a different type of document, a
3 document describing a dataset, and asked for that dataset as
4 used in a model, and it was not. I -- that's my memory is
5 there's one or two examples of that.

6 THE COURT: Uh-huh.

7 MR. TUTTLE: But if the Plaintiffs chose a dataset
8 described on --

9 THE COURT: A data card.

10 MR. TUTTLE: -- a data card for a model, then that
11 dataset is for that model, and we will -- and we will have
12 found it, and we will have provided it.

13 THE COURT: Okay.

14 MR. TUTTLE: And if there -- if anything comes up,
15 if there's any issue, we will disclose that to Plaintiffs,
16 and we will work through it. But I -- the only ones I'm --
17 the only ones I can think of right now where I think we have
18 not found that dataset in connection to that -- with that
19 model was where the Plaintiffs pointed to a document that
20 was not describing a particular model, right? It was
21 describing a dataset or a project or something, and they
22 said they wanted that dataset, and it -- that was not used
23 in -- this is my memory, it was not used in connection with
24 that model.

25 THE COURT: Okay.

1 MR. TUTTLE: And, again, I think the Court is
2 quite right to focus on the class definition that the Court
3 has -- that Judge Lee specifically has given the Plaintiffs'
4 leave to amend to and once the motion to dismiss is ruled
5 on, which is US persons who own a US copyright in any work
6 used by Google to train Google's generative AI models. So
7 what we're focused on are the works that were used to train,
8 and the training datasets for the models will tell us, by
9 what's in them, what those works are. And the source code
10 does not tell us that.

11 THE COURT: Okay. Mr. Young made the point about
12 the mitigation filters.

13 MR. TUTTLE: Yes, your Honor.

14 THE COURT: And that the source code has the
15 instructions as to when and how to use those. So what --
16 does the training -- how do the -- it doesn't sound like the
17 datasets or the data cards would contain that information.

18 MR. TUTTLE: Well, the purpose of the mitigation
19 -- that mitigation software is to use the training data that
20 was used for the models in order to create -- in order to be
21 able to detect, you know, repetition or duplication of the
22 training data for that model in the output of that model.
23 So the code for these mitigation measures will contain the
24 detailed algorithms for how an index would be built, how it
25 would prevent recitation in the output of a model, but it

1 won't tell you -- it won't tell you the data, and the data
2 -- and, again, the data is the training data for the model.
3 That's the whole point, right, is to have the -- is to have
4 those -- that mitigation software able to detect the
5 presence of the same data that went into the model, right?
6 Because it's trying to -- it's trying to block, you know,
7 repetition of that data in the output.

8 THE COURT: So if I'm understanding your
9 explanation, it's you have the content of the training data.

10 MR. TUTTLE: Yes.

11 THE COURT: That's what we've been talking about.
12 You have the exact content. To the extent that there's a
13 mitigation filter to limit the output from that training
14 data, that's okay. But the training data is what it is.

15 MR. TUTTLE: That's correct, your Honor. And I --
16 you know, I just -- the source code isn't going -- again,
17 all -- the source code will get into the gory, technical
18 details of how things are done, but it will not, at the end
19 of the day, tell us what works, are involved, or anything
20 like that.

21 THE COURT: Uh-huh. Uh-huh. Uh-huh. Okay. I
22 started you off by answering some of my questions. Did you
23 have a further response to Mr. Young's presentation?

24 MR. TUTTLE: Well, I don't want to -- I would
25 prefer to focus on the things, if anything, that concerned

1 the Court or had questions. I don't want to -- I don't want
2 to spend time on issues that aren't important. I will just
3 say, like, walking through the history of the crawling --
4 again, the source code for the crawling won't -- as the
5 declaration said, it won't tell us what URLs were crawled.
6 I don't think the Plaintiffs are asking for all the URLs
7 that Google has ever crawled in the history of the internet.
8 That's what Google does. Google tries to crawl as much of
9 the internet as possible. That's how search works. But in
10 any event, it doesn't relate to the class definition. The
11 class definition isn't a class of all people whose
12 copyrighted works were crawled by Google. If they -- it's
13 all people whose works were used to train Google's
14 generative AI models. So there's all these steps, right?

15 THE COURT: Uh-huh.

16 MR. TUTTLE: But at the end of the day, the
17 training data that is used to train the model is the source
18 for the works that are used for that purpose.

19 THE COURT: Mr. Young pointed to what he -- may be
20 an inconsistency in the Zhao declaration, and that's docket
21 152-2, specifically with regards to the language in
22 paragraph three. And I'm quoting from the third sentence of
23 paragraph three,

24 "These crawlers handle everything from
25 domain resolution and protocol

1 negotiation to parsing the robots.text
2 protocol, extracting the URLs and
3 content from web pages, and crawl
4 scheduling and rate limiting."

5 And as I understood Mr. Young's argument, he was
6 focused on the extracting URLs language and acknowledging in
7 paragraph five that the declarant says that specific URLs of
8 websites is not demonstrated in the source code.

9 So can you explain -- can those positions be harmonized
10 or can you explain what -- is it Mr. Zhao or Doctor Zhao is
11 referring to in paragraph three with regards to the
12 extraction of URL?

13 MR. TUTTLE: Yes, I think that Doctor Zhao is
14 explaining that the source code for the crawlers will
15 contain the routines that do the extraction of the URL, but
16 they don't -- but the URLs are not themselves in the code.
17 So the -- when a Google crawler crawls the web, right, it
18 extracts the content of websites and metadata about the
19 websites, and then that information is stored, and -- but
20 the code --

21 THE COURT: Is the instruction to do it.

22 MR. TUTTLE: -- is the instruction on how to do
23 that extraction. It doesn't contain the extracted text or
24 metadata, which -- such as the URL.

25 And then what matters at the end of the day is what

1 training data got selected and used to train the models,
2 right? And so that's, again, I think why I keep focusing
3 and the Court seems focused on that training data --

4 THE COURT: Uh-huh.

5 MR. TUTTLE: -- which is sort of the end of the
6 line of what is -- of what the models are trained on.

7 THE COURT: Provided there's metadata that shows
8 where it came from.

9 MR. TUTTLE: I understand, your Honor, and I
10 believe there will be. If it turns out that there's some
11 example where that's not present, then we will look to see
12 where -- if there are -- if there's other information that's
13 available that would contain it, but it would not be in the
14 source code regardless.

15 THE COURT: Uh-huh.

16 MR. TUTTLE: So the source code is not -- I
17 believe it will be in the training data --

18 THE COURT: I understand.

19 MR. TUTTLE: -- but if it turns out not to be,
20 we'll have to figure out some other -- if there's something
21 else. But in no event --

22 THE COURT: Are you aware of anything else? I
23 mean, I know you said, "Well, we don't have a" -- I think
24 you said words to the effect of, "There's not a list of
25 every URL that crawl or data has ever visited," which if its

1 job is to go out into the internet everywhere, that would --
2 I understand that. But are you aware of any other tracking
3 -- other than the metadata in the training sets, which we --
4 one hopes, but may or may not be there, then is there some
5 other source that you're aware of?

6 MR. TUTTLE: I am not. The only -- I mean, the --
7 well, we would -- and we would have to investigate. There
8 could be a -- like a -- an earlier version, like an earlier
9 extraction or something like that that could contain it. If
10 that exists, we could look for that.

11 THE COURT: Earlier than what? Than the training
12 dataset?

13 MR. TUTTLE: Correct.

14 THE COURT: Uh-huh.

15 MR. TUTTLE: If it's -- but, again, the -- that
16 won't -- that's not the code. It --

17 THE COURT: I understand.

18 MR. TUTTLE: -- it would be an earlier -- like a
19 super set of the training data that existed at an earlier
20 stage, if that still exists. That would be an option, and
21 we could look for other ways if this comes up. I hope that
22 it will not.

23 THE COURT: I hope so too.

24 MR. TUTTLE: But it won't be in the source code
25 regardless.

1 THE COURT: Okay. All right. That is helpful.

2 Let me get some final comments from Mr. Young, and

3 I'll --

4 MR. TUTTLE: Thank you. Oh -- I'll -- no. That's
5 it, your Honor.

6 THE COURT: Okay.

7 MR. YOUNG: Thank you, your Honor. Christopher
8 Young on behalf of Plaintiffs. And I just have a couple of
9 comments as I was kind of just sitting here and reflecting
10 on the colloquy. One of the things I just want to make sure
11 we recall -- we all collectively remember is that class
12 membership is not the only thing we will have to show at
13 class certification. We'll have to show that common issues
14 predominate over individual ones. And it is our contention
15 that the source code would be the best evidence of that,
16 because the source code is really the best reflection of the
17 conduct at issue, which is not just the copying, but the
18 continual use. And I think under the case law we cited in
19 our discovery brief, it is recognized that each use must be
20 justified, and the source code really is the best evidence
21 of the conduct of each infringing use, or perhaps each use
22 that Google would say would be fair use. That's the first
23 point.

24 THE COURT: But the source code is -- I mean, I
25 always think of it as a -- not being a coder by nature, but

1 as a set of instructions. I mean, it's -- you know, it sets
2 parameters, direction, and then it goes out and executes.

3 MR. YOUNG: Certainly.

4 THE COURT: I mean, it's executing an algorithm,
5 really, right?

6 MR. YOUNG: Certainly. But that's exactly why it
7 would be relevant for class cert, understanding what that
8 common course of conduct is, especially when the documents
9 themselves may not evidence some of the things that will be
10 important for class cert. And by way of example, your
11 Honor, I'll point you to -- you know, pending right now in
12 San Francisco right now in federal court is Judge Alsup is
13 wrestling with class certification right now in a case
14 involving Anthropic. And, you know, he --

15 THE COURT: I am aware.

16 MR. YOUNG: Yeah. And one of the -- one of the --
17 one of the back and forth between the Court and the parties
18 between the filings is whether or not there should be a
19 pirate -- pirated works class and a non-pirated works class.
20 And if -- what the source code will show and what we've --
21 at least I have not seen from the documents is any
22 acknowledgment of whether or not Google has taken any
23 pirated information at all. But as Mr. Tuttle has -- I
24 think no one disagrees that Google's entire -- one of
25 Google's main businesses is to scrape the internet. But if

1 they're scraping pirate websites with intent and knowledge
2 that they're doing so to use for training the models, that
3 could be -- that could be a dividing line for certain
4 classes in this case.

5 THE COURT: And I appreciate that, and we -- I
6 know your papers address the subclass and need for
7 information, but if the dataset has metadata -- if the
8 training data has the metadata, then you'll have that
9 information.

10 MR. YOUNG: Potentially.

11 THE COURT: Yeah, ideally.

12 MR. YOUNG: I mean, it depends where we get the
13 training dataset. I mean, it was unclear to me now whether
14 we will have kind of the training data as it was fed before
15 the model, after it gets cleaned up because -- you know,
16 just again, it's taking a step back from the life cycle of,
17 kind of, data, right? You have data as it's scraped, right?
18 And I imagine there's going to be some cleanup. I mean, I
19 do not think, for example, Google is feeding, like, URLs --
20 constantly feeding URLs into the model to train, because --
21 I just don't think from, like, a practical matter. That'll
22 just make the thing essentially just a URL regurgitation
23 device, and that's not how we understand these things to do,
24 right? So if the -- I mean, I just do not know at what
25 stage the training datasets are going to be produced to us

1 in a form that is -- that we're going to be inspecting. For
2 -- so for us, which is why we're also asking for the
3 crawlers, which at our first instance will tell us kind of
4 the provenance of it. And we --

5 THE COURT: I understand.

6 MR. YOUNG: Yeah.

7 THE COURT: I understand. Okay.

8 MS. WEAVER: Just on the issue of the pirated and
9 the copyrighted, if there is source code that shows
10 filtering after the training set comes, that is very
11 significant for the case because it tells you intent, it
12 tells you that Google can filter copyrighted material out
13 and chose not to do it on the front end, and it also gets at
14 there's copyrighted material, and then there's known pirated
15 material. And the URLs alone won't show that, and it won't
16 get at the heart of whether they're distinguishing it at any
17 point in the process.

18 It also gets at the transformative nature, you know,
19 what are they doing with the data? What is the Court doing?
20 We did have a proffer when we were meeting and conferring
21 from Google to say they would not challenge this
22 transformative element of class cert. And if we work that
23 out and put that in a stipulation and order for the Court,
24 maybe that could bear on that discrete issue. But there is
25 certainly a bucket of source code that is giving

1 instructions about the use of categories that really go to
2 the heart of the case.

3 MR. YOUNG: And, your Honor, just to put a finer
4 point on what my colleague, Ms. Weaver, just said, you know,
5 it is our understanding, based on our meet and confer
6 earlier today, that Google is not going to challenge at
7 class cert that -- the individual -- for each individual
8 models, what happens kind of in -- at the model level, the
9 neural network level, the conversion from training data to
10 output would not be an individualized issue, which is why
11 we're currently dropping, at least for class certification
12 purposes, our ask for the code, the evidence, kind of the
13 training and data processing. So that kind of intermediate
14 step in between the collection, collating, and out to the
15 last step of kind of the mitigation.

16 THE COURT: Okay. I appreciate that. All right.
17 Anything further, Mr. Young?

18 MR. YOUNG: That's it, your Honor. Thank you.

19 THE COURT: Okay. All right.

20 All right. Thank you all very much for that. Let me
21 take a short recess just to organize my notes and see if I
22 have any further follow-up questions while I have everyone
23 here. So we'll take 10 minutes and -- that might be 15, but
24 I'm going to shoot for 10.

25 And Ms. (Indiscernible), if you'll come back and get me

1 then, I'll appreciate it.

2 All right. Thank you. We're in recess. You may
3 remain seated.

4 (Proceedings recessed briefly.)

5 THE COURT: Great. Thank you for your patience.
6 I wanted a few minutes to go back through my notes and the
7 materials that the parties had submitted in support of their
8 respective positions with regards to the production of
9 source code at this stage.

10 I do think that the review process that's been
11 established, there are undoubtedly some wrinkles in that,
12 there always are in these inspection protocols, but it seems
13 like at least it is set up to address a number of the
14 concerns raised by Plaintiffs as to insight into the data
15 that they need, certainly for purposes of class
16 certification. It may not be complete, that is the datasets
17 able to be accessed through this environment, but it's too
18 early to tell. It looks like the objective is for it to be
19 complete, in that the datasets -- the actual content is
20 provided and the metadata for that content.

21 So with that, I don't see yet the necessity for source
22 code, at least in the -- certainly not in the scope as the
23 current request is pending. So I'm going to deny the
24 request for source code without prejudice -- I want --
25 provided that Plaintiffs' access to the datasets in the

1 environment that's been described here today by Google is
2 available next week. And I know it's going to take some
3 time to get oriented and start to look at that and evaluate
4 what is there, and I'm -- obviously, it's in Plaintiffs'
5 best interest to jump on that and move through that. And if
6 there are data points that go to issues of class
7 certification that are missing or cannot be discerned, you
8 need to meet and confer on those. But this has to be
9 prompt. We're all on the clock here. And response to that
10 has to be prompt. And if there is, again, something that
11 can be identified or a piece of necessary data for class
12 cert that can be identified that's not available through
13 this environment, inspection of datasets, and -- then
14 Plaintiffs can come back with a request for the source code,
15 and it should -- it would be more limited and focused on
16 what is missing. But I'm hopeful that through the meet and
17 confer process, the parties could come to some agreement.

18 It may well be that at some point in this process that
19 source code has to be produced and made available. It's not
20 unusual in these cases. But let's get a look at what of the
21 needs are met through the dataset inspection as it is set
22 up. So that'll get started next week.

23 If Plaintiffs have identified a dataset that Google
24 cannot confirm was used in training -- and I understand your
25 point, Mr. Tuttle, that to the extent datasets are picked

1 off of data cards, they were used in training. But if they
2 come from some other source, if Google cannot confirm that
3 they were used in training, then that needs to be
4 communicated promptly here at the start of the review
5 process, and the Plaintiffs can replace that dataset with
6 some other selection, because obviously we need -- you need
7 the information that was actually used in training.

8 All right. So I am going to issue a -- I'll issue a
9 summary order from today that captures these points with
10 regards to both source code and the custodian issue, so
11 everyone's record is clear and the docket is clear. And --
12 but my rulings are as I have articulated them from the bench
13 here today.

14 All right. Any questions? Anything that I've
15 just said that you're thinking, "Oh, no, that can't possibly
16 happen." And hearing none, excellent. Thank you.

17 I want to, again, thank Counsel for your preparation,
18 both in the papers and the supporting materials for today's
19 positions -- presentations, as well as the presentations
20 themselves, and for restating your names every time you were
21 at the podium for the benefit of my clerks.

22 All right. Thank you all very much. Have a good day.

23 MR. YOUNG: Thank you, your Honor.

24 MS. WEAVER: Thank you, your Honor.

25 MR. TUTTLE: Thank you, your Honor.

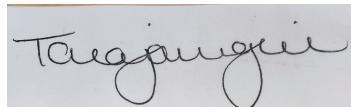
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

THE CLERK: We are adjourned.
(Proceedings adjourned at 11:49 a.m.)

CERTIFICATE OF TRANSCRIBER

I certify that the foregoing is a true and correct transcript, to the best of my ability, of the above pages of the official electronic sound recording provided to me by the U.S. District Court, Northern District of California, of the proceedings taken on the date and time previously stated in the above matter.

I further certify that I am neither counsel for, related to, nor employed by any of the parties to the action in which this hearing was taken; and, further, that I am not financially nor otherwise interested in the outcome of the action.

A handwritten signature in cursive script, appearing to read "Teagunzie", is centered within a light gray rectangular box.

Echo Reporting, Inc., Transcriber

Monday, June 23, 2025